

# Embodiment of Learning in Electro-Optical Signal Processors

Michiel Hermans,<sup>1,\*</sup> Piotr Antonik,<sup>1,†</sup> Marc Haelterman,<sup>2</sup> and Serge Massar<sup>1</sup>

<sup>1</sup>*Laboratoire d'Information Quantique, Université Libre de Bruxelles,  
50 Avenue F. D. Roosevelt, CP 225, B-1050 Brussels, Belgium*

<sup>2</sup>*Service OPERA-Photonique, Université Libre de Bruxelles,  
50 Avenue F. D. Roosevelt, CP 194/5, B-1050 Brussels, Belgium*

Delay-coupled electro-optical systems have received much attention for their dynamical properties and their potential use in signal processing. In particular it has recently been demonstrated, using the artificial intelligence algorithm known as reservoir computing, that photonic implementations of such systems solve complex tasks such as speech recognition. Here we show how the backpropagation algorithm can be physically implemented on the same electro-optical delay-coupled architecture used for computation with only minor changes to the original design. We find that, compared when the backpropagation algorithm is not used, the error rate of the resulting computing device, evaluated on three benchmark tasks, decreases considerably. This demonstrates that electro-optical analog computers can embody a large part of their own training process, allowing them to be applied to new, more difficult tasks.

Published version : <http://dx.doi.org/10.1103/PhysRevLett.117.128301>

*Introduction.* Nonlinear dynamical systems, such as neural networks (NN), can be used to perform highly complex computations, e.g. speech or image recognition. One of the main difficulties when using such systems is to train their internal parameters. The backpropagation (BP) algorithm [1, 2] is one of the most important algorithms in this area, and is behind the remarkable successes achieved in the field of deep learning in the last decade [3]. The simple idea behind the BP algorithm is to compute the derivative (or gradient) of a cost function in the parameter space of the system. The gradient is then subtracted from the parameters themselves in order to reduce the cost function. This process is repeated until the cost function no longer reduces.

Such nonlinear dynamical systems can be implemented in hardware. Here also the training of internal parameters is key and the use of the BP algorithm is highly beneficial in order to improve performance [4, 5]. However implementing the BP algorithm in hardware systems can be difficult because of the need of an accurate model to compute the gradient and because of the resources necessary to run the BP algorithm. Remarkably, in certain cases the BP algorithm can be implemented physically on the system it is optimising [6]. The basic idea behind this advance is to use a slightly modified version of the system for propagating error signals backwards, i.e. for running the BP algorithm. Such self-learning computing systems could be highly advantageous, as any gain in terms of processing speed or limited power consumption will also apply to the training phase. Furthermore having the same hardware computing the BP algorithm eliminates, to a large extent, the need for an accurate model of the system. This idea may conceivably also have implications for biological neural networks, as these are physical system that – using mechanisms that are not yet well understood – can both compute and carry out

their own training process. Reference [6] also reported a proof of concept experiment in which physical BP was tested on a simple task, but left open the question of whether the algorithm, with all the imperfections inherent in an experiment, can provide the same improvement in performance as numerical approaches [4, 5].

References [4–6] used as computational device a delay dynamical system (see [7, 8]). Such systems can be exploited to realise a form of analog computer based on the Reservoir Computing (RC) paradigm [9, 10] in which unoptimised high-dimensional dynamic systems (termed *reservoirs*) are used as signal processors. The RC approach is simple, versatile and can be applied to a wide set of problems (see the review [11]) and experimental implementations [12–20]. Applying the BP algorithm to delay-coupled signal processors allows one to optimise many more parameters than in traditional RC, yielding significant improvements in performance as was shown in simulation in [4], and subsequently in an experiment [5] in which BP was applied to a numerical model of the system, and the results of the BP algorithm applied to the physical experimental setup.

Here we implement the BP algorithm physically on an electro-optic delay dynamical system used as signal processor. Our key innovation is to modify the system used in [16, 17] by adding a photonic setup capable of implementing both the nonlinearity and its derivative, so that it can be used both as signal processor and to perform the BP algorithm. We test our system on several tasks considered hard in the machine learning community, including a real world phoneme recognition task (the TIMIT task, discussed later in this paper), obtaining state of the art results when the BP algorithm is used. The present work thus demonstrates the full potential of physical BP. It constitutes an important step towards self-learning hardware, with potential applica-

tions towards ultra-fast, low energy consumption, computing systems.

In the following we first recall the principles of reservoir computing and error back propagation, before introducing our experimental implementation. We then report the results obtained on several benchmark tasks, and conclude with a discussion of the results and their implications.

*Reservoir Computing.* In typical RC tasks, the goal is to map an input sequence  $s_i$  (where  $i \in \{1, \dots, L\}$ , with  $L$  the total sequence length) to an output sequence  $y_i$ , which has target values  $y_i^*$ , for example a speech signal to a sequence of labels. In order to use delay-coupled systems as reservoir computers, the discrete time input sequence  $s_i$  is encoded into a continuous time function  $z(t)$  by the input mask  $m(r)$  and bias mask  $m_b(r)$ , where  $r \in [0, T]$ , with  $T$  the *masking period*, as follows

$$z(t) = z(iT + r) = m(r)s_i + m_b(r) . \quad (1)$$

In our implementation, we use a delay-coupled system with sine nonlinearity (which stems from the transfer function of the intensity modulator, as will be explained below), which obeys the equation:

$$a(t + D) = \mu \sin(a(t) + z(t)) \quad (2)$$

where  $a(t)$  is the state variable and  $D$  is the delay. The factor  $\mu$  corresponds to the total loop amplification. Eq. (2) can be seen as a special case of the Ikeda delay differential equation [21].

One then needs to map the continuous time state variable  $a(t)$  to a discrete time output sequence  $y_i$ . This is performed using an output mask  $u(r)$  where  $r \in [0, T]$  and a bias term  $u_b$  as follows:

$$y_i = \int_0^T dr a(iT + r)u(r) + u_b . \quad (3)$$

In the RC paradigm the input mask is typically chosen randomly, and the output mask  $u(r)$  and  $u_b$  is determined by solving a linear system of equations which minimises the mean square error  $C$  between the desired and actual output:  $C = \langle (y_i - y_i^*)^2 \rangle_i$ .

*Error Backpropagation.* The goal of applying error backpropagation to the above scheme is to optimise both the input and output masks  $m(r)$ ,  $m_b(r)$ ,  $u(r)$  and  $u_b$ , knowing the output  $a(t)$ , and the desired output  $y_i^*$ . To this end one needs the gradient of  $C$  with respect to the masks, given by (the proof is given in the Supplementary

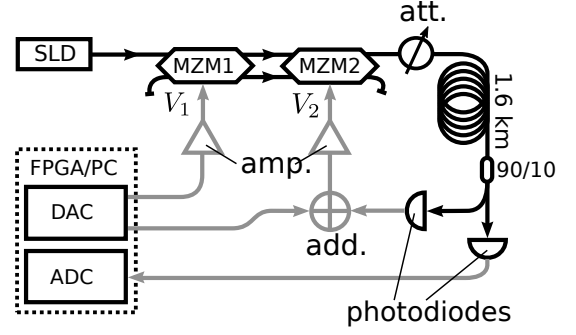


FIG. 1: Schematic representation of the experimental system. SLD: superluminescent diode; MZM1 and MZM2: dual input - dual output Mach-Zehnder Modulators;  $V_1$  and  $V_2$ : driving voltages of the MZMs; att.: programmable optical attenuator; add.: electrical combiner; amp.: pulse amplifier.

Material):

$$\bar{e}(iT + r) = e_i u(r) , \quad (4)$$

$$e(t - D) = J(t) (e(t) + \bar{e}(t)) , \quad (5)$$

$$J(t) = \mu \cos(a(t) + z(t)) , \quad (6)$$

$$\frac{dC}{dm(r)} = \sum_i e(iT + r)s_i , \quad (7)$$

$$\frac{dC}{dm_b(r)} = \sum_i e(iT + r) , \quad (8)$$

where  $\bar{e}(t) = \partial C / \partial a(t)$  is a continuous time signal and, as above,  $i \in \{1, \dots, L\}$  and  $r \in [0, T]$ . One can then iteratively improve the masks so as to lower  $C$ .

*Physical BP.* In order to use the same hardware for both the signal processing and its own training, one exploits the very close analogy between Eqs. (1) and (4) – both are formed in the same way from a discrete time sequence, multiplied by a periodic mask – as well as the very close analogy between Eqs. (2) and (5) – both are delay systems. However the equation for  $e(t)$  depends on future values, so it needs to be solved backwards in time. In practice one time-inverts  $\bar{e}(t)$  and  $J(t)$  before computing  $e(t)$  to obtain a linear delayed equation:

$$e(q + D) = J(q) (e(q) + \bar{e}(q)) . \quad (9)$$

where we use  $q$  instead of  $t$  to remind oneself that we are dealing with time-inverted signals. We also note that  $J(t)$ , the derivative of the nonlinear function, is a cosine, which can also be implemented using the intensity modulator. Although this property of the sine function is key for this experiment, other types of nonlinearity can be implemented in analogue hardware (see the discussions).

*Experimental implementation.* In the present work we show how Eqs. (2) and (9) can be realised using the same physical setup. Our fibre optics experiment is depicted in Figure 1. Light is generated by a superluminescent

diode (SLD) emitting in the telecommunications band (1550 nm, with a 33 nm FWHM), modulated by two dual input / dual output Mach-Zehnder modulators (MZM), and attenuated using a programmable optical attenuator used to control the total loop amplification of the system, i.e.  $\mu$  in Eq. (2). It then propagates through an approximately 1.6 km long spool of optical fibre which provides a total loop delay of 7.93  $\mu$ s. The light is split and enters two photodiodes, one of which provides the feedback signal. The signals are produced and recorded by Digital to Analog Converters (DAC) and Analog to Digital Converters (ADC), controlled by a Xilinx Virtex 6 FPGA chip. The FPGA simultaneously generates the input voltage signals and records the output signals. The FPGA communicates with a PC that controls the whole experiment. (Further details on the experimental setup are given in the Supplementary material).

The key innovation with respect to the earlier experiments [16, 17] is the use of two dual input / dual output MZMs, see Fig. 1, which allows to implement both Eqs. (2) and (9) using the same physical system. Taking into account the incoherence of light in the two branches between the modulators (see Supplementary Material for details), the output of the upper branch of MZM2 (see Fig. 1) can be found to be:

$$I_2^+ = \frac{I_0}{2} [1 + \sin(V_1/V_0) \sin(V_2/V_0)], \quad (10)$$

where  $I_0$  is the input intensity in the upper branch of MZM1,  $V_1$  and  $V_2$  are the driving voltages and  $V_0$  a constant depending on the MZM. The computational details are presented in the Supplementary Material. In the forward mode, we choose  $V_1/V_0 = \pi/2$ . The transfer function thus acts as a sinusoidal function for the input argument  $V_2/V_0 = a(t) + z(t)$ . The constant offset  $I_0/2$  is removed by the high-pass filter of the amplifier, that drives the MZM. Therefore, once the loop is closed, we end up with Eq. (2). In the backward mode we drive MZM1 with a voltage  $V_1/V_0 = a(q) + z(q) + \pi/2$ , and MZM2 with a signal proportional to  $\bar{e}(q) + e(q)$ , but scaled down sufficiently such that  $\sin(V_2/V_0) \approx V_2/V_0 = \bar{e}(q) + e(q)$ , which gives the desired functionality for the adjoint system Eq. (9).

In order to train our reservoir computer, we first choose a value of  $\mu$  close to the threshold for instability. We then iterate the following three steps for (typically) several thousands of iterations, during which performance slowly improves until it converges:

- 1) We take the training data (typically a small subsequence of the complete set), and convert it to  $z(t)$  using the input masks. We feed this signal to the experimental setup, physically implementing Equation 2. Next, we measure and record the signal  $a(t)$ , and generate an output sequence  $y_i$  using the output masks.

- 2) From the output and the desired target values we compute the sequence  $e_i = \partial C / \partial y_i$  at the output, and

convert it to  $\bar{e}(t)$ , now using the output mask as an input mask. Next we time-invert it and feed it back into the experimental setup. Simultaneously we drive the first MZM with the (time-inverted) signal  $a(q) + z(q)$  in order to implement the online multiplication with  $J(q)$ . We record the response signal  $e(q)$ .

- 3) From the recorded signals  $a(t)$  and  $e(t)$  we obtain the gradients for the masking signals, which we use to update the input and output masks:

$$\begin{aligned} m(r) &\leftarrow m(r) - \eta dC/dm(r), \\ m_b(r) &\leftarrow m_b(r) - \eta dC/dm_b(r), \\ u(r) &\leftarrow u(r) - \eta dC/du(r), \\ u_b &\leftarrow u_b - \eta dC/du_b, \end{aligned} \quad (11)$$

where  $\eta$  is a (typically small) learning rate. In order to speed up convergence we applied a slightly more advanced variant of these update rules known as Nesterov momentum [22, 23] (details are given in the Supplementary material).

*Results.* We experimentally validate the above scheme using the system described in Fig. 1 by testing it on three time series processing task. We consider first of all the NARMA10 task [24], an academic task often used in the RC community. Here the input sequence  $s_i$  consists of a series of independent and identically distributed random numbers drawn uniformly from the interval  $[0, 0.5]$ . The desired output sequence is given by

$$y_i^* = 0.3y_{i-1}^* + 0.05y_{i-1}^* \sum_{n=1}^{10} y_{i-n}^* + 1.5s_i s_{i-9} + 0.1.$$

The second task we will call VARDEL5 (from *variable delay*). Here the input sequence consists of i.i.d. digits drawn from the set  $\{1, 2, 3, 4, 5\}$ . The desired output is then given by  $y_i^* = s_{i-s_i}$ , i.e., the goal is to retrieve the input instance delayed with the number of time steps given by the current input.

As a performance metric for NARMA10 and VARDEL5 we use the *normalised root mean square error* (NRMSE), which is given by

$$\text{NRMSE} = \sqrt{\frac{\langle (y_i - y_i^*)^2 \rangle_i}{\langle (y_i^*)^2 \rangle_i}}.$$

The NRMSE varies between 0 (perfect match), and 1 (no relation between output and target).

The third task is a frame-wise phoneme labelling task. We use the TIMIT dataset [25], a speech dataset in which each time step has been labelled with one of 39 phonemes. The input data is high-dimensional (consisting of 39 frequency channels), and the desired output is one of (coincidentally) 39 possible output classes. The goal is to label each frame in a separate test set. Consequently, the performance metric is now the classification error rate, i.e., the fraction of misclassified phonemes in the test

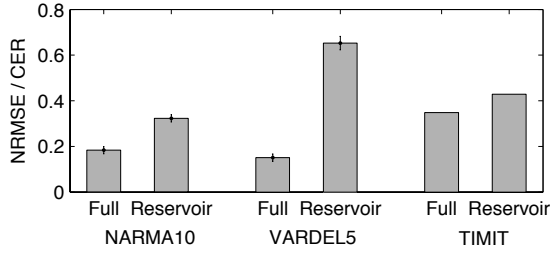


FIG. 2: Comparison of performances for the three tasks under consideration. We show either NRMSE (for NARMA10 and VARDEL5) or the classification error rate (CER) for TIMIT. For each task we show performance for a fully trained systems (Full) vs. those trained using the RC paradigm (Reservoir). Error bars indicate standard deviations if available.

set. Note that the masking scheme and BP algorithm is easily extended to multidimensional in – and output sequences (more details are provided in the Supplementary material). The TIMIT task has been studied before in the context of RC, which has shown it to be challenging, typically requiring extremely large reservoirs to obtain competitive performance [26, 27].

For all these tasks we compared performance of the fully trained system to traditional RC, where we kept the input and bias masks fixed and random, and only optimised their global scaling and the feedback strength parameter  $\mu$ . Full experimental details for each task may be found in the supplementary material, together with an example of optimised input masks and of convergence of NRMSE during training. The results are shown in Figure 2. The experimental setup is successful in performing both useful computations, and implementing its own training process. The fully trained system consistently outperforms the RC approach in all tasks considered.

For the NARMA10 task we improve over all previous experimental results. The previous best was published in [20], which reported an NRMSE of 0.249 for 50 virtual nodes, and 0.22 for 300 virtual nodes, whereas here we obtain a NRMSE of 0.185 for 80 nodes (note that in [20] they report *normalised mean square error* (NMSE), which is the square of the NRMSE). That result was obtained on an experimental setup that was specially designed to produce a minimal amount of noise (using a passive cavity as a reservoir). The lowest reported experimental NRMSE on a setup equivalent to ours was 0.41 [17]. Note that we obtain a better average performance for the RC setup (NRMSE = 0.32), which is most likely due to the higher number of virtual nodes (80 as opposed to 50 in [17]).

For the VARDEL5 task, we cannot directly compare to literature, however as pointed out in chapter 5 of [28], this task is an important example of a task that is so nonlinear that it is nearly impossible to solve it with RC.

This is confirmed here; the NRMSE of RC is 0.66, indicating that the reservoir has only captured the task on a very rudimentary level. The fully trained system shows a drastically better performance (NRMSE = 0.15). This shows that training the input masks not just allows for better performance on existing tasks, but also allows to tackle tasks that are so intricate that they are considered beyond the reach of traditional RC.

For the TIMIT task we obtain a classification error rate of 34.8% for fully trained systems, vs. 42.9% for the standard RC approach. These results are only slightly worse than similar experimental results presented in [5], (33.2% for fully trained systems and 40.5% for the RC approach) where 600 virtual nodes were used as opposed as 200 in our case.

*Discussion.* The present work confirms the results anticipated in [4, 5]: the performance of delay-based reservoir computers can be drastically improved by optimising both input and output masks. Furthermore, following the proposal of [6], we showed that the underlying hardware is capable of running a large part of its own optimisation process. We performed our demonstrations on a fast electro-optical system (whose speed could be readily improved by several orders of magnitude, see, e.g. [15]), and on tasks considered hard in the RC community. Importantly, our work has revealed that the BP algorithm is robust against various experimental imperfections (see the Supplementary Material for details), as the performance gains we obtained on all three tasks were similar to those predicted by numerical simulations.

Although our experiment relies on the sine nonlinearity and its cosine derivative, other nonlinear functions can also be successfully realised in hardware with their derivatives. For instance, the so-called linear rectifier function, which truncates the input signal below a certain threshold, is a popular activation function in neural architectures [29]. Its derivative is a simple binary function which can be easily implemented using an analogue switch, as in [6]. In [30] it is shown how to implement a sigmoid nonlinearity and its derivative. And in [18, 20] the nonlinearity is quadratic, and therefore the derivative, which is linear, should also be easy to implement. Furthermore, the BP algorithm is robust against imperfect implementation of the derivative, as shown in section 4.3 of the Supplementary Material, and in the Supplementary Material of [6] (Supplementary Note 4). Therefore we expect that physical implementation of the BP algorithm will be possible in a wide variety of physical systems.

The current setup still requires some slow digital processing to perform the masking and to compute gradients from the recorded signals. Performing masking operations in analog hardware, however, is actively being researched [31], and these approaches could be used to speed up the present setup. Another limitation is the relatively slow data transfer between the FPGA and the



computer. Implementing the full training algorithm on the FPGA would drastically increase the speed of the experiment. FPGA's have already been demonstrated to be useful for controlling and training electro-optical signal processors [32, 33].

Nowadays, there's an increased interest in unconventional, neuromorphic computing, as this could allow for energy efficient computing, and may provide a solution to the predicted end of Moore's law [34]. These novel approaches to computing will likely be made with components that exhibit strong element-to-element variability, or whose characteristics evolve slowly with time. Self-learning hardware may be the solution that enables these systems to fulfil their potential. The results in [6] and in this paper therefore constitute an important step towards this goal.

The authors acknowledge financial support by Interuniversity Attraction Poles Program (Belgian Science Policy) project Photonics@be IAP P7-35, by the Fonds de la Recherche Scientifique FRS-FNRS and by the Action de Recherche Concertée of the Fédération Wallonie-Bruxelles through grant AUWB-2012-12/17-ULB9.

---

\* michiel.hermans@ulb.ac.be

† piotr.antonik@ulb.ac.be

- [1] D. Rumelhart, G. Hinton, and R. Williams, *Learning internal representations by error propagation* (MIT Press, Cambridge, MA, 1986).
- [2] P. Werbos, *Neural Networks* **1**, 339 (1988).
- [3] Y. LeCun, Y. Bengio, and G. Hinton, *Nature* **521**, 436 (2015).
- [4] M. Hermans, J. Dambre, and P. Bienstman, *Neural Networks and Learning Systems*, *IEEE Transactions on* **26**, 1545 (2015), ISSN 2162-237X.
- [5] M. Hermans, M. C. Soriano, J. Dambre, P. Bienstman, and I. Fischer, *JMLR* **16**, 2081 (2015).
- [6] M. Hermans, M. Burm, T. Van Vaerenbergh, J. Dambre, and P. Bienstman, *Nature communications* **6**, 6729 (2015).
- [7] T. Erneux, *Applied delay differential equations* (Springer Science & Business Media, 2009).
- [8] V. Flunkert, I. Fischer, and E. Schöll, *Philosophical Transactions of the Royal Society of London A: Mathematical, Physical and Engineering Sciences* **371**, 20120465 (2013).
- [9] H. Jaeger, Tech. Rep. GMD Report 152, German National Research Center for Information Technology (2001).
- [10] D. Verstraeten, B. Schrauwen, M. d'Haene, and D. Stroobandt, *Neural Networks* **20**, 391 (2007).
- [11] M. Lukoševičius and H. Jäger, *Computer Science Review* **3**, 127 (2009).
- [12] C. Fernando and S. Sojakka, in *Proceedings of the 7th European Conference on Artificial Life* (2003), pp. 588–597.
- [13] K. Caluwaerts and B. Schrauwen, in *Proceedings of the 2nd International Conference on Morphological Computation* (2011).
- [14] L. Appeltant, M. C. Soriano, G. Van der Sande, J. Danckaert, S. Massar, J. Dambre, B. Schrauwen, C. R. Mirasso, and I. Fischer, *Nature Communications* **2**, 468 (2011).
- [15] D. Brunner, M. C. Soriano, C. R. Mirasso, and I. Fischer, *Nature Communications* **4**, 1364 (2013).
- [16] L. Larger, M. Soriano, D. Brunner, L. Appeltant, J. Gutierrez, L. Pesquera, C. Mirasso, and I. Fischer, *Optics express* **3**, 20 (2012).
- [17] Y. Paquot, F. Duport, A. Smerieri, J. Dambre, B. Schrauwen, M. Haelterman, and S. Massar, *Scientific Reports* **2**, 1 (2012).
- [18] K. Vandoorne, P. Mechet, T. Van Vaerenbergh, M. Fiers, G. Morthier, D. Verstraeten, B. Schrauwen, J. Dambre, and P. Bienstman, *Nature Communications* **5** (2014).
- [19] N. D. Haynes, M. C. Soriano, D. P. Rosin, I. Fischer, and D. J. Gauthier, *Physical Review E* **91**, 020801 (2015).
- [20] Q. Vinckier, F. Duport, A. Smerieri, K. Vandoorne, P. Bienstman, M. Haelterman, and S. Massar, *Optica* **2**, 438 (2015).
- [21] K. Ikeda and K. Matsumoto, *Physica D: Nonlinear Phenomena* **29**, 223 (1987).
- [22] Y. Nesterov, *Soviet Mathematics Doklady* **27**, 372 (1983).
- [23] I. Sutskever, J. Martens, G. Dahl, and G. Hinton, in *Proceedings of the 30th international conference on machine learning (ICML-13)* (2013), pp. 1139–1147.
- [24] A. Atiya and A. Parlos, *IEEE Transactions on Neural Networks* **11**, 697 (2000).
- [25] J. Garofolo, N. I. of Standards, T. (US, L. D. Consortium, I. Science, T. Office, U. States, and D. A. R. P. Agency, *TIMIT Acoustic-phonetic Continuous Speech Corpus*. (Linguistic Data Consortium, 1993).
- [26] F. Triefenbach, A. Jalalvand, B. Schrauwen, and J.-P. Martens, in *Advances in Neural Information Processing Systems 23* (2010), pp. 2307–2315.
- [27] F. Triefenbach, K. Demuyne, and J.-P. Martens, *IEEE Signal Processing Letters* **21**, 311 (2014).
- [28] M. Hermans, Ph.D. thesis, Ghent University (2012), URL <http://hdl.handle.net/1854/LU-3171075>.
- [29] X. Glorot, A. Bordes, and Y. Bengio, in *14th International Conference on Artificial Intelligence and Statistics* (2011), vol. 15, p. 275.
- [30] B. Shi and C. Lu, *Generator of neuron transfer function and its derivative* (2002), US Patent 6429699.
- [31] F. Duport, A. Smerieri, A. Akrou, M. Haelterman, and S. Massar, *Scientific Reports* **6** (2016).
- [32] P. Antonik, F. Duport, A. Smerieri, M. Hermans, M. Haelterman, and S. Massar, in *APNNA's 22th International Conference on Neural Information Processing* (2015), vol. 9490 of *LNCIS*, pp. 233–240.
- [33] P. Antonik, M. Hermans, F. Duport, M. Haelterman, and S. Massar, in *SPIE's 2016 Laser Technology and Industrial Laser Conference* (2016), vol. 9732.
- [34] M. M. Waldrop, *Nature* **530**, 144147 (2016).

# Embodiment of Learning in Electro-Optical Signal Processors. Supplementary Material.

Michiel Hermans, Piotr Antonik, Marc Haelterman, Serge Massar

October 28, 2016

## 1 Experimental setup

The experimental setup depicted in Fig. 1 of the main text uses the following components:

- Superluminescent diode (SLED): Thorlabs, model SLD1550P-A40), center wavelength 1550 nm, FWHM 33 nm.
- Dual input/dual output Mach Zehnder modulators (MZM): EOspace, model number AX-2x2-0MSS-12-PFA-PFA.
- Programmable optical attenuator: Agilent, model 81571A.
- Photodiodes: Terahertz Technologies, model TIA-525.
- Pulse amplifiers: Mini-Circuits, model ZHL-32A+.
- FPGA, ADC, and DAC: 4DSP FMC151 daughter card containing a two-channel DAC and ADC, controlled by a Xilinx Virtex 6 FPGA chip.

Note that the FPGA simultaneously generates the voltage signal that represents  $z(t)$  and records the voltage signal representing  $a(t)$ . The FPGA also performs a minimal signal processing step by selecting and averaging over the middle samples of each masking step (see Section 4.1 for more details). The remaining processing steps are carried out on a PC.

Sending and receiving data to and from the FPGA is currently the main speed bottleneck of the experiment. Even though a single training iteration lasts only about 0.6 seconds for the NARMA10 and VARDEL5 task, most of this time is spent on the communication overhead with the PC (buffering). If the entire experiment were to be performed on the FPGA (which is feasible), a single training iteration would take of the order of milliseconds.

## 2 Online multiplication using cascaded MZMs

As mentioned in the main text, we use two back to back dual input/dual output Mach-Zehnder modulators for implementation of both Eqs. (2) and (10) from the main text using the same setup. The main fact we use is that the spectrum of the SLD is narrow enough to allow for a large extinction ratio by the MZMs, but is broad enough that

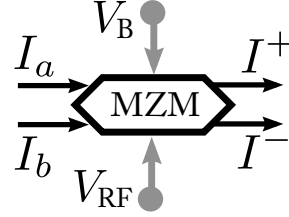


Figure 1: Schematic representation of a dual input / dual output Mach Zehnder modulator. The MZM is driven by the sum of two input voltages: one constant bias voltage  $V_B$  and a fast signal  $V_{RF}$ .

the light in the two branches entering MZM2 from MZM1 can be considered incoherent. In the present experiment, the coherence length of the light from the SLD is of the order of a few hundreds of micrometers, which means that a very small difference in path length for the connections in between MZM1 and MZM2 is sufficient to make the two signals incoherent.

Consider the operation of a single MZM, schematised in Fig. 1. The intensities of the incoming light sources are denoted by  $I_a$  and  $I_b$ , and the MZM is driven by a voltage  $V$ , which is the sum of a constant bias voltage  $V_B$  and a fast voltage signal  $V_{RF}$ . The bias voltage was omitted in the main text to avoid confusion. Taking into account the incoherence between the two input signals, the intensities of the output branches ( $I^+$  and  $I^-$ ) are given by:

$$\begin{aligned} I^+ &= I_a \frac{1 + \sin(V/V_0)}{2} + I_b \frac{1 - \sin(V/V_0)}{2}, \\ I^- &= I_a \frac{1 - \sin(V/V_0)}{2} + I_b \frac{1 + \sin(V/V_0)}{2}, \end{aligned} \quad (1)$$

with  $V_0$  a constant depending on the MZM.

It is now easy to model the output of the two cascaded MZMs. Suppose the source has an intensity  $I_0$ , and no light enters the second input of MZM1. And suppose MZM1 and MZM2 receive voltages  $V_1$  and  $V_2$ , respectively. The output intensities  $I_1^+$  and  $I_1^-$  of MZM1 are given by

$$\begin{aligned} I_1^+ &= I_0 \frac{1 + \sin(V_1/V_0)}{2}, \\ I_1^- &= I_0 \frac{1 - \sin(V_1/V_0)}{2}. \end{aligned} \quad (2)$$

The intensity  $I_2^+$  at the first output branch of MZM2 is then:

$$I_2^+ = I_0 \frac{(1 + \sin(V_1/V_0))(1 + \sin(V_2/V_0))}{4} \quad (3)$$

$$+ I_0 \frac{(1 - \sin(V_1/V_0))(1 - \sin(V_2/V_0))}{4} \quad (4)$$

$$= \frac{I_0}{2} [1 + \sin(V_1/V_0) \sin(V_2/V_0)]. \quad (5)$$

In the experiment MZM1 receives a constant bias signal on top of an RF driving signal, such that  $V_1/V_0 = \pi/2 + V'_1/V_0$ , with  $V'_1$  the RF signal. We can thus write:

$$I_2^+ = \frac{I_0}{2} [1 + \cos(V'_1/V_0) \sin(V_2/V_0)].$$

We use the setup in two modes. In the forward mode,  $V'_1 = 0$ , so that the cascaded MZMs behave as:

$$I_2^+ = \frac{I_0}{2} [1 + \sin(V_2/V_0)],$$

i.e., the transfer function acts as a sinusoidal function for the input argument  $V_2/V_0$ , which is equal to the sum of the input signal  $z(t)$  and the system state  $a(t)$ . Note that a constant offset  $I_0/2$  is added to the output. We use, however, amplifiers with a high-pass filter to drive the MZMs, which remove the DC offset. Therefore, once the loop is closed, this constant bias is removed, and we effectively end up with Eq. (2) in the main text.

In the backwards mode, we drive MZM1 with a voltage  $V'_1$  proportional to  $a(q - D) + z(q - D)$ . MZM2 is driven with a signal proportional to  $\bar{e}(q) + e(q)$ , but scaled down sufficiently such that  $\sin(V_2/V_0) \approx V_2/V_0 = \bar{e}(q) + e(q)$ . This means that in the backwards mode we can write:

$$I_2^+ = \frac{I_0}{2} [1 + \cos(a(q + D) + z(q + D)) (\bar{e}(q) + e(q))],$$

which is (up to the constant bias, and the factor  $\mu$  which is imposed later by the optical attenuator) the desired functionality for the adjoint system (see Eq. (10) in the main text).

### 3 Derivation of gradients and adjoint system

#### 3.1 Setting up the problem

We wish to find the gradient of a cost function  $C$  w.r.t. the parameters that can be optimised. In order to achieve this we have to use the chain rule through all the dependencies that describe the system. We will then obtain the backward equations given in the main text. Figure 2 gives a schematic of how the forward and backward equations must be implemented experimentally. Figure 3 depicts the information flow in the forward and backward systems.

We first recall the relevant equations describing the forward system. The input signal  $z(t)$ , formed by concatenating the input masks weighted with the current input sample  $s_i$  can be rewritten as

$$z(t) = s_{\lceil t/T \rceil} m(t \bmod T) + m_b(t \bmod T), \quad (6)$$

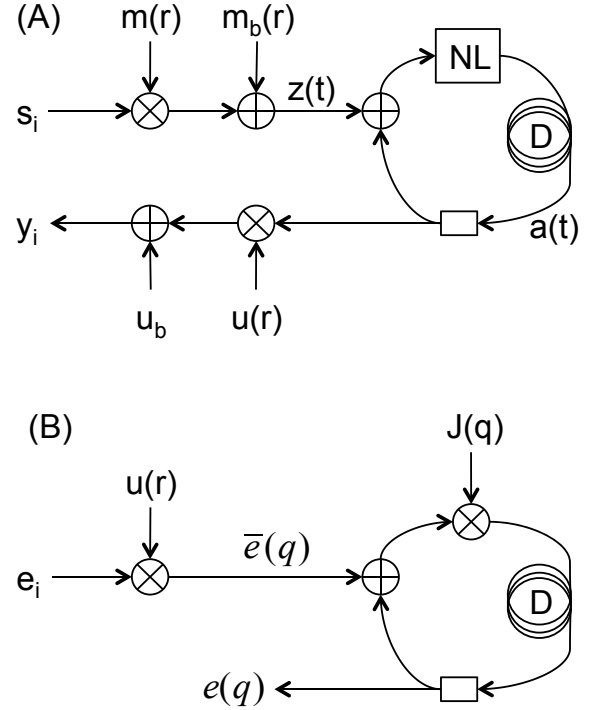


Figure 2: **A:** Schematic depiction of the forward system, as given in the main text and in Supplementary Material Eqs. (6, 7, 8). **B:** Schematic depiction of the backward system, as given in the main text and in Supplementary Material Eqs. (12, 17), where  $q$  is the backwards time.

where  $\lceil \cdot \rceil$  indicates the ceiling function, so that  $\lceil t/T \rceil = i$  gives the index of  $s_i$  corresponding to the time  $t$ . We use the modulo operation in the argument of the input masks to indicate that the masks are repeated over time. Next we write down the expression for the reservoir state  $a(t)$ :

$$a(t + D) = \mu \sin(a(t) + z(t)). \quad (7)$$

Finally, we can write the formula for the output instances  $y_i$  as follows:

$$y_i = u_b + \int_0^T dr u(r) a_i(r), \quad (8)$$

with  $a_i(r) = a(r + (i-1)T)$ , the  $i$ -th segment of the recording of  $a(t)$ .

In what follows, for the sake of generality and of simplicity of notation, we take the input and output masks to be continuous functions of time. We denote functional derivatives with respect to time dependent functions as ordinary derivatives. The case, relevant to practical implementations, in which the masks depend on a finite number of parameters, is discussed in section 4.1. For simplicity in the derivations we will assume, unless indicated otherwise, that all variables, both in continuous time  $t$  and discrete time  $i$ , are defined for  $i$  and  $t$  going from  $-\infty$  to  $\infty$ . If we have a specific finite input sequence  $s_i$  with  $i \in \{1, \dots, L\}$ , we simply extend this beyond these bounds assuming that all extra  $s_i$  are equal to zero. Similarly, we assume that  $z(t)$  is zero if

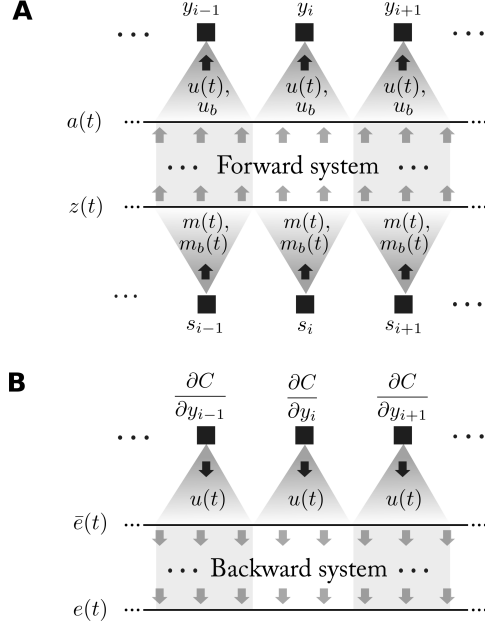


Figure 3: **A:** Schematic depiction of information flow when the system is used in the forward direction. On the bottom, the input sequence  $s_i$  is converted to a continuous-time signal  $z(t)$  (with time running from left to right). Each instance in the sequence is multiplied with the finite-length masking signal  $m(t)$  and added to  $m_b(t)$ . These sequences are then concatenated in time to form  $z(t)$ , the input to the forward system. The output  $a(t)$  of the forward system is then converted into an output sequence  $y_i$  by segmenting  $a(t)$  in time, and multiplying the segments with the output masks  $u(t)$ , and integrating over each of them. **B:** Schematic depiction of the information flow in the “backwards” mode. The derivatives  $\partial C/\partial y_i$  are used as an input sequence. They are multiplied by  $u(t)$  which now plays the role of input mask. This yields the signal  $\bar{e}(t)$  that serves as input for the backward system. The output of the backward system is  $e(t)$ .

$\lceil t/T \rceil \notin \{1, \dots, L\}$ . Subsequently, if we sum or integrate over  $i$  or  $t$  without indicating limits, this indicates a summation or integration from  $-\infty$  to  $\infty$ . In practice it turns out that if we only have a finite sequence, we only need to compute states over its corresponding time span. Similarly, when performing backpropagation, we only need to compute backwards over the same time span. All states outside of this interval do not influence the gradient computation, which means there are no problems in considering only finite intervals. This matters as in realistic training scenarios we typically train on relatively short sequences (in the case of the present paper of length 100).

### 3.2 Output mask gradient

For the output masks we can write

$$\frac{dC}{du(r)} = \sum_i \frac{\partial C}{\partial y_i} \frac{dy_i}{du(r)}. \quad (9)$$

For example, if the cost function we wish to minimise is the squared error over the interval of the input sequence:

$$C = \sum_{i=1}^L (y_i - y_i^*)^2,$$

$$e_i = \frac{\partial C}{\partial y_i} = 2(y_i - y_i^*) \text{ for } i \in \{1, \dots, L\}.$$

$$\frac{\partial C}{\partial y_i} = 0 \text{ for } i \notin \{1, \dots, L\}.$$

The second factor in Eq. (9) we can get from Equation 8:

$$\frac{dy_i}{du(r)} = a_i(r),$$

such that the gradient for the output mask  $u(t)$  is simply given by

$$\frac{dC}{du(r)} = \sum_i \frac{\partial C}{\partial y_i} a_i(r),$$

or, given the fact that  $\partial C/\partial y_i = 0$  outside the interval in which the sequence is defined:

$$\frac{dC}{du(r)} = \sum_{i=1}^L \frac{\partial C}{\partial y_i} a_i(r). \quad (10)$$

Similarly we find that

$$\frac{dC}{du_b} = \sum_{i=1}^L \frac{\partial C}{\partial y_i}.$$

### 3.3 Input mask gradient

The case of the input masks is more involved. Working out the chain rule we find:

$$\begin{aligned} \frac{dC}{dm(r)} &= \sum_i \frac{\partial C}{\partial y_i} \frac{dy_i}{dm(r)} \\ &= \sum_i \frac{\partial C}{\partial y_i} \int dt' \frac{\partial y_i}{\partial a(t')} \frac{da(t')}{dm(r)} \\ &= \int dt' \bar{e}(t') \frac{da(t')}{dm(r)}, \end{aligned} \quad (11)$$

where we have used

$$\bar{e}(t') = \frac{\partial C}{\partial a(t')} = \sum_i \frac{\partial C}{\partial y_i} \frac{\partial y_i}{\partial a(t')}.$$

From Equation 8 we can obtain (using a modulo function in the argument of  $u(r)$ ):

$$\frac{\partial y_i}{\partial a(t')} = \delta_{i, \lceil t'/T \rceil} u(t' \bmod T),$$

i.e., equal to zero when  $t'$  did not fall in the segment of time used to produce  $y_i$ , and equal to the output mask otherwise. This yields:

$$\begin{aligned} \bar{e}(t') &= u(t' \bmod T) \frac{\partial C}{\partial y_{\lceil t'/T \rceil}} \\ &= u(r) e_i \end{aligned} \quad (12)$$



where  $r = t' \bmod T$  and  $i = \lceil t'/T \rceil$ . In other words,  $\bar{e}(t)$  is produced by masking the sequence  $\partial C/\partial y_i$  with the output mask  $u(r)$ .

The second factor in Equation 11 we work out as follows. Using the chain rule we get

$$\frac{da(t')}{dm(t)} = \int dt'' \frac{da(t')}{dz(t'')} \frac{dz(t'')}{dm(t)}. \quad (13)$$

and

$$\frac{da(t')}{dz(t'')} = \frac{\partial a(t')}{\partial a(t'')} + \int dt''' \frac{\partial a(t')}{\partial a(t''')} \frac{da(t''')}{dz(t'')}.$$

From equation 7 we obtain the partial derivatives:

$$\begin{aligned} \frac{\partial a(t')}{\partial z(t'')} &= \frac{\partial a(t')}{\partial a(t'')} \\ &= \mu \delta(t' - t'' - D) \cos(a(t' - D) + z(t' - D)), \end{aligned}$$

Or, more compactly:

$$\frac{\partial a(t')}{\partial z(t'')} = \delta(t' - t'' - D) J(t'),$$

with

$$J(t') = \mu \cos(a(t' - D) + z(t' - D)).$$

This yields

$$\frac{da(t')}{dz(t'')} = J(t') \left[ \delta(t' - t'' - D) + \frac{da(t' - D)}{dz(t'')} \right].$$

By filling in the expression for  $da(t' - D)/dz(t'')$  recursively we can write this as:

$$\frac{da(t')}{dz(t'')} = \sum_{i=0}^{\infty} \left[ \delta(t' - t'' - iD) \prod_{j=0}^{i-1} J(t' - jD) \right]. \quad (14)$$

By filling in Equation 14 in Equation 13, and inserting the result in Equation 11 we obtain:

$$\frac{dC}{dm(r)} = \int dt' dt'' \bar{e}(t') \sum_{i=0}^{\infty} \delta(t' - t'' - iD) \prod_{j=0}^{i-1} J(t' - jD) \frac{dz(t'')}{dm(r)}. \quad (15)$$

We can solve the integral over  $t'$  explicitly. We denote this by  $e(t'')$ :

$$\begin{aligned} e(t'') &= \int dt' \bar{e}(t') \sum_{i=0}^{\infty} \delta(t' - t'' - iD) \prod_{j=0}^{i-1} J(t' - jD) \\ &= \sum_{i=0}^{\infty} \bar{e}(t'' + iD) \prod_{j=0}^{i-1} J(t'' + (i - j)D) \\ &= \sum_{i=0}^{\infty} \bar{e}(t'' + iD) \prod_{j=1}^i J(t'' + jD). \end{aligned} \quad (16)$$

It's straightforward to prove that  $e(t)$  is equal to the expression as presented in the main text (with arguments shifted by  $D$ ):

$$e(t) = J(t + D)(e(t + D) + \bar{e}(t + D)). \quad (17)$$

Indeed, if we recursively fill in the expression for  $e(t + D)$  in Eq. 17, we obtain Eq. 16. Using this we can reduce Equation 15 to

$$\frac{dC}{dm(r)} = \int dt'' e(t'') \frac{dz(t'')}{dm(r)}.$$

From the expression of  $z(t)$  we find that

$$\frac{dz(t'')}{dm(r)} = \delta(t'' \bmod T - r) s_{\lceil t''/T \rceil}.$$

Inserting this we can find the final expression for the gradient for the input mask:

$$\frac{dC}{dm(r)} = \sum_i s_i e_i(r),$$

or, again using the fact that we defined  $s_i = 0$  for  $i \notin \{1, \dots, L\}$ :

$$\frac{dC}{dm(r)} = \sum_{i=1}^L s_i e_i(r), \quad (18)$$

with  $e_i(r) = e(r - (i - 1)T)$ , the  $i$ -th segment of the time trace of  $e(t)$ . Similarly for  $m_0(t)$  we can write:

$$\frac{dC}{dm_b(r)} = \sum_{i=1}^L e_i(r). \quad (19)$$

### 3.4 Multiple inputs/outputs

The above explanation is easily extended to multiple input and output dimensions. Suppose we have a multivariate time series  $\mathbf{s}_i$ , where the  $k$ -th element at time step  $i$  is denoted by  $\mathbf{s}_i[k]$ . We can then easily construct  $z(t)$  by defining as many input masks  $m_k(t)$  as there are input dimensions and adding them all up:

$$z(t) = \sum_k \mathbf{s}_{\lceil t/T \rceil}[k] m_k(t \bmod T) + m_b(t \bmod T),$$

The desired output can similarly exist of a multivariate time series with elements  $\mathbf{y}_i^*[l]$ . To produce an output  $\mathbf{y}_i[l]$  we simply define an output mask  $u_l(t)$  and bias  $u_l^0$  for each output channel:

$$\mathbf{y}_i[l] = u_l^0 + \int_0^T dt u_l(r) a_i(r).$$

The same procedure can now be used to determine the gradients with respect to the multivariate input and output masks. We find:

$$\frac{dC}{du_l(r)} = \sum_{i=1}^L \frac{dC}{d\mathbf{y}_i[l]} a_i(r), \quad (20)$$

and

$$\frac{dC}{du_l^0} = \sum_{i=1}^L \frac{dC}{d\mathbf{y}_i[l]}. \quad (21)$$

The source of the BP equation is now

$$\bar{e}(t') = \sum_l u_l(t' \bmod T) \frac{dC}{dy_{\lceil t'/T \rceil}[l]}, \quad (22)$$

the recurrence for the error  $e(t)$ , eq. (17), is unchanged, and one has

$$\frac{dC}{dm_k(r)} = \sum_{i=1}^L \mathbf{s}_i[k] e_i(r).$$

## 4 Implementation details

### 4.1 Mask parametrisation

While the aforementioned theory is generally valid for continuous-time signals, an experimental setup is limited by the finite bandwidth of the DAC/ADC, and the analog electronic parts. To make sure that these effects play a limited role, we parametrise the input and output masks as piecewise constant functions, which has been common practice for reservoirs of this type [3]. To this end we divide the delay  $D$  into an integer number  $N_D$  of equal time segments, called *masking steps*. Next we ensure that the masking period  $T$  has a total duration that is also contains an integer number  $N_T$  of masking steps, in our case one less than the delay:  $N_T = N_D - 1$ . This allows for the mixing of the states over time, as detailed in [3].

The input and output masks are picked to be constant for the duration of each masking step. This implies that  $z(t)$  is piecewise constant. The fact that both  $T$  and  $D$  are an integer number of masking steps makes that changes in  $a(t)$  only occur in between the masking steps, i.e., they are synchronised with the masking steps, and this is valid for the backwards pass too. In short,  $a(t)$ ,  $\bar{e}(t)$  and  $e(t)$  are all piecewise constant signals, with values that remain constant during each masking step.

In practice this allows us to reduce effects of noise by averaging the signals representing  $a(t)$  and  $e(t)$  over several measuring samples during a single masking step. Typically we pick a set of samples from the middle of each masking step, and discard those at the beginning and the end as they may contain artefacts caused by the limited bandwidth of the ADC. More importantly, it allows us to make a discrete time approximation of the entire system. For example, let's consider equation 8. The mask  $u(t)$  is made up of  $N_T$  constant segments of equal length, with values during the segments denoted  $u_k$ . Similarly, each segment  $a_i(r) = a(t - (i - 1)T)$  is piecewise constant, with values we can for example denote with  $a_k^i$ . The integral reduces to

$$y_i = u_b + \sum_{k=1}^{N_T} a_k^i u_k,$$

(where we absorbed the factor  $T$  that emerges from the integration into the values  $u_k$ ). Each particular value  $a_k^i$  can be interpreted as the state of the  $k$ -th ‘neuron’ or ‘node’ state during the  $i$ -th instance of the input sequence. We can still use the expressions for the gradients in Equations 10, 18 and 19. Indeed, by construction, the gradient for the

output mask  $u(t)$  for the duration of a single masking step is a constant (as  $a(t)$  remains constant over the segment). The same holds for the gradients for the input masks. This implies that  $u(t)$  and  $m(t)$  remain piecewise constant during training, and we can in practice describe them simply as lists of values instead of a continuous-time function.

Note that the choice of dealing with bandwidth limitations by using piecewise constant functions is not the only possible avenue. One alternative would be to impose bandwidth constraints on  $m(t)$  and  $u(t)$ , such that the finite signal generator bandwidth and sampling rates form no obstacle in treating the setup as a continuous-time setup. We chose the piecewise-constant constraint as it is more directly related to existing implementations of delay-coupled electro-optical signal processors, and it allows to identify a specific number of ‘virtual nodes’ (the number of segments within the masking period  $T$ ). In other words, the choice of  $N_T$  determines the ‘complexity’, or the number of degrees of freedom of the system.

### 4.2 Gradient descent

We used stochastic gradient descent to train the masks; each iteration we drew a 100 time step sequence to determine a gradient. This sequence was either generated on the fly (in the case of VARDEL5 and NARMA10), or drawn randomly from a training set (TIMIT). Note that as the BP equation is linear, we are in principle free to rescale  $\bar{e}$  as we wish. In practice, in order to keep MZM2 in the linear regime, we scaled the input error signal  $\bar{e}(t)$  by dividing it by its standard deviation and multiplying with a factor 0.1. The learning rate  $\eta$  we choose equal to 0.25 at the start of the training process, after which it drops linearly to zero throughout the course of the experiment. On top of that we use Nesterov momentum with a momentum factor 0.9 to speed up convergence [2, 4]. Nesterov momentum is a heuristic method that finds widespread use in speeding up convergence of stochastic gradient descent. The idea of momentum in gradient descent is to give parameter updates a certain inertia, meaning that previous parameter updates still count in the current one, which helps with overcoming local minima and speeds up convergence. Nesterov momentum is a simple variation of this principle, where the algorithm measures the gradient one update step ahead in order to change its momentum “ahead of time”.

### 4.3 Robustness

Our work shows that physical BP is robust against imperfections of the physical setup, as illustrated by the following imperfections we were confronted with.

The first imperfection was the high-pass filtering operation of the amplifiers used to drive the MZMs, with a cut-off frequency of 20 kHz. While the high-pass filter is a desirable property (to get rid of voltage bias), this corresponds to a typical time scale of about 8  $\mu$ s, which is about the same as the loop delay and therefore not negligible. The current experimental setup does not take this filtering operation into account explicitly.

A second imperfection was an imbalance in losses between the two fibres connecting MZM1 with MZM2. A third imperfection was that the system was not perfectly linear during the backwards pass, since MZM2 is never a perfectly linear system. There's also an important trade-off here. One can reduce the residual nonlinearity by reducing the amplitude of the incoming voltage signal that represents  $\bar{e}(t)$ . But in turn this also reduces the signal-to-noise ratio of the measurement during the backpropagation phase, such that one needs to find a good balance between these two effects.

All these effects are imperfections inherent to the physically implemented backpropagation phase, but both in simulation and in the actual experiments we found that they only had a very minor impact on the training process and the overall performance.

One parameter that turned out to be crucial was the bias voltage of MZM2. The reason is that even a small offset from an effectively zero level introduces a systematic error in the backpropagation process, such that the measured signal (denoted as  $e_c(t)$  to indicate that it is corrupted) becomes :

$$e_c(t - D) = J(t)(e_c(t) + \bar{e}(t) + \tilde{e}),$$

with  $\tilde{e}$  a constant offset caused by an incorrectly set voltage bias of MZM2. It turned out that, in the experiments, keeping this bias level effectively equal to zero was difficult; very slight drifts on the effective working point of the MZM occurred over the course of minutes/hours. Luckily, the backpropagation is a linear process. This means that we can recover  $e(t)$  by performing a second measurement right after measuring  $e_c(t)$ :

$$e_r(t - D) = J(t)(e_r(t) + \tilde{e}),$$

and

$$e(t) = e_c(t) - e_r(t).$$

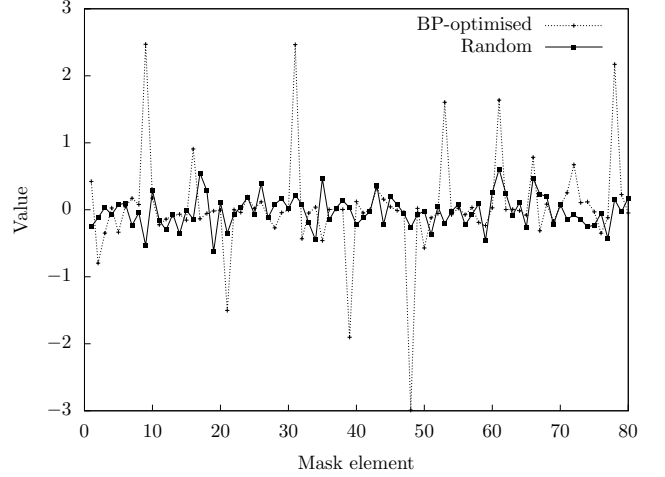
In other words we simply need to perform two measurements after each other, where in the second one we send a 'zero' input error, and subtract this from the first measurement in order to remove the influence of the offset of MZM2. This turned out to solve the problem.

## 5 Tasks

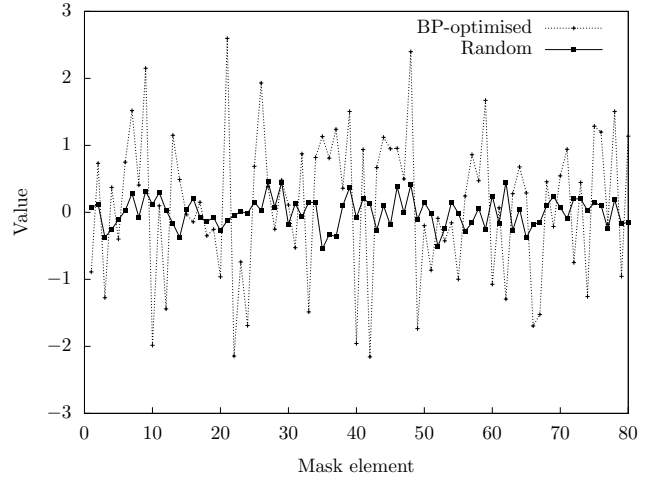
### 5.1 NARMA10 and VARDEL5

In the case of NARMA10 and VARDEL5 we divided  $T$  into 80 equal time intervals ( $N_T = 80$ ), which allowed us to take 16 samples during each masking step, where we averaged over the middle 8 in order to get piecewise-constant values for  $a(t)$  and  $e(t)$ . We chose the number of training iterations at 10,000, 20,000 for VARDEL5 and NARMA10, respectively, chosen heuristically as a trade-off between the time required for an experiment and the final performance. (A single iteration lasted approximately 0,6 s).

The cost functions used for NARMA10 and VARDEL5 are the aforementioned sums of squared errors. We repeated the training cycles 10 times, each time with different random input mask initialisations. Output masks were always



(a) Input mask  $m$



(b) Bias input mask  $m_b$

Figure 4: Comparison of BP-optimised input masks (dotted curves) and random RC masks (solid curves) for the VARDEL5 task, for  $m(r)$  (top panel) and  $m_b(r)$  (bottom panel).

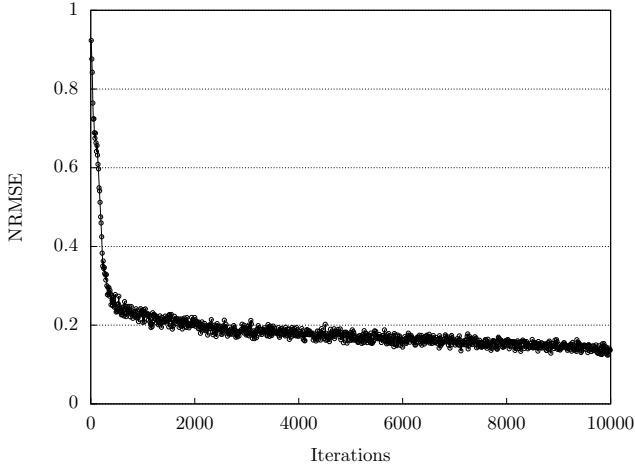


Figure 5: Evolution of the NRMSE during the training process on the VARDEL task. The error falls sharply during the first 300 iterations, and then converges slowly towards 0.15.

initialised at zero. For all backpropagation experiments we set the feedback parameter strength parameter  $\mu$  effectively equal to one (such that the system is at the ‘edge of stability’), which we found to give the best performance.

In Figure 4 we depict the masks  $m(r)$  and  $m_b(r)$  for the RC implementation (when they are chosen at random), and after optimisation using the BP algorithm, for the VARDEL5 task. One sees that the BP algorithm dramatically changes the input masks. In particular the mask  $m$  is very large at some specific values of  $r$ , and almost zero for other values. This suggests that in some sense what the optimised reservoir is doing is storing the value of the input on specific neurons, and then keeping it in memory for some time, before mixing it nonlinearly with the input several time steps in the future. In Figure 5 we depict how the NRMSE converges over time for the VARDEL5 task, as the BP algorithm slowly improves the input and output masks.

For the reservoir computing results we measured average performance as a function of three scaling parameters: the feedback strength parameter  $\mu$  and the scaling of the input mask and bias mask. Once optimal parameters were determined we ensured that the output masks were trained on an unlimited amount of input training data (in practice we observed the test error for increasing amounts of training data, and stopped as soon as the performance no longer improved). This was to ensure that we have a fair comparison to the backpropagation setup, where we generate unlimited amounts of data too. Each experiment was repeated 10 times, giving rise to the error bars in Figure 3A in the main text.

## 5.2 TIMIT

For the TIMIT task we used 1,000,000 training iterations. Because of this large number of iterations we only performed a single full training cycle.

Measurement noise plays a smaller role in a classification task such as this one, and we divided  $T$  into 200 masking

steps, taking 8 samples in each and averaging over the middle 4, thereby increasing the number of virtual nodes  $N_T$  while taking into account the hardware constraints (sample rates of the DAC and ADC). Most likely this number can be increased further for example using only 4 samples per masking steps and averaging over the middle 2. In practice we are also limited by the relatively slow communication between the PC and the FPGA, and increasing  $N_T$  increases the amount of data that needs to be transferred, slowing down the experiment considerably. Currently, for 1,000,000 iterations the training took two weeks to complete.

We picked  $\mu$  at a value slightly under one, but we found in simulations that performance did not strongly depend on it for a broad range of values.

In the case of TIMIT, the goal is to minimise a classification error rate, which is not directly differentiable. One possible strategy is simply to try and minimise the MSE between the output and the target labels (1 for the correct class, zero for all others). Classification would then be performed by the *winner-take-all* approach, where we simply select the output channel with the highest output as the ‘winner’. In practice, using MSE for classification suffers from some drawbacks. Most importantly MSE will put a lot of emphasis on producing the exact target values (close to zero or one), while we are only interested in performance after selecting the highest output. A better approach is to use a softmax function at the output, which converts the output values into a set of probabilities, and minimise the cross-entropy with the target probabilities (again, 1 for the correct class and zero for all others). Details on this strategy can be found for example in [1]. In practice the conversion of the output  $\mathbf{y}_i[k]$  into probabilities is performed using the so-called *softmax* function:

$$\mathbf{p}_i[k] = \frac{\exp(\mathbf{y}_i[k])}{\sum_l \exp(\mathbf{y}_i[l])},$$

The cost function is the cross-entropy:

$$C = - \sum_{i=1}^L \sum_k \mathbf{t}_i[k] \ln \mathbf{p}_i[k],$$

where we denote the target outputs as  $\mathbf{t}_i[k]$ . It can then be shown that

$$\frac{dC}{d\mathbf{y}_i[k]} = \mathbf{p}_i[k] - \mathbf{t}_i[k],$$

(and again zero if  $i \notin \{1, \dots, L\}$ ) This means that the error we have at the output takes on virtually the same form as before, only this time there is the intermediary step of the softmax function. Gradients for the output masks are almost the same as before, except for equations 20 and 21 where we use  $\mathbf{p}_i[k] - \mathbf{t}_i[k]$  instead of  $\mathbf{y}_i[k] - \mathbf{y}_i^*[k]$ . As far as the rest of the BP algorithm goes, we now simply have to mask these ‘output errors’ to produce  $\bar{e}(q)$ , and the rest plays out exactly as before.

For the RC approach, optimising the parameters (input scaling, bias scaling and feedback gain) on the hardware would be too costly in terms of time. Therefore we optimised them on a PC using a simulation of the physical



setup. Once we decided on the parameters, we ran all the TIMIT data through the physical setup and recorded all the responses. Next we trained output weights, again using gradient descent with the above cross-entropy loss.

## References

- [1] Christopher M. Bishop, *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer, August 2006.
- [2] Yurii Nesterov. *A method of solving a convex programming problem with convergence rate  $o(1/k^2)$* . Soviet Mathematics Doklady, 27(2) : 372–376, 1983.
- [3] Yvan Paquot, Francois Duport, Antoneo Smerieri, Joni Dambre, Benjamin Schrauwen, Marc Haelterman, and Serge Massar. *Optoelectronic reservoir computing*. Scientific Reports, 2: 1–6, 2012.
- [4] Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. *On the importance of initialization and momentum in deep learning*. In Proceedings of the 30th international conference on machine learning (ICML-13), pages 1139–1147, 2013.